



Bacterial strain Taxonomy for Genomic Surveillance

Enterobase and Hierarchical Clustering

Nigel Dyer: University of Warwick, UK

Created: October 2025

Outline



This session consists of the following elements

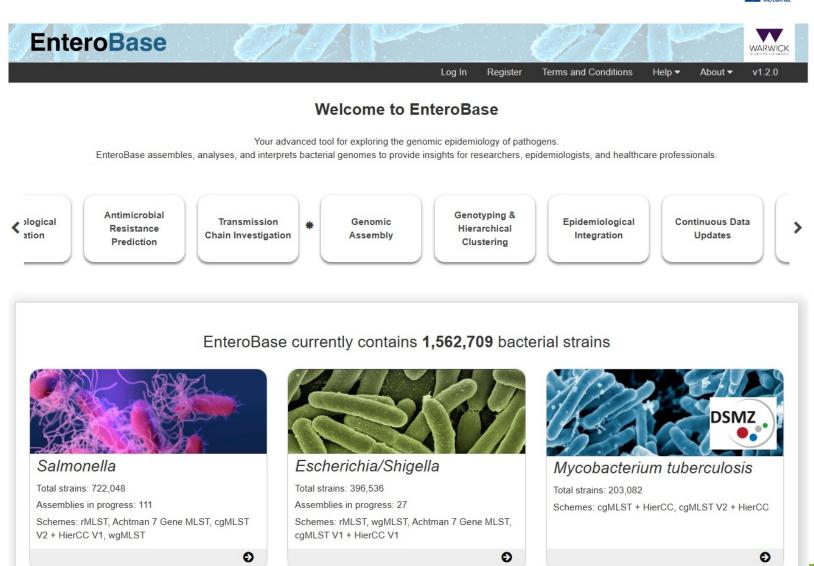
- 1. Introduction to Enterobase
- 2. Hierarchical clustering and LIN codes
- 3. Usage examples

EnteroBase — https://enterobase.warwick.ac.uk/



Species:

- Salmonella
- Escherichia
- Streptococcus
- Clostridioides
- Vibrio
- Heliobacter
- Yersinia
- Moraxella







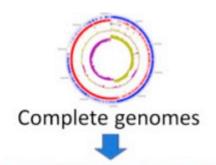
Are you familiar with EnteroBase?









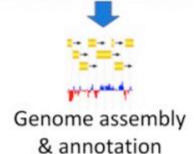


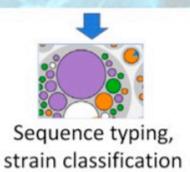


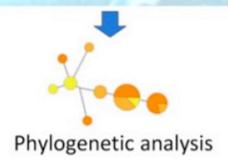
EnteroBase

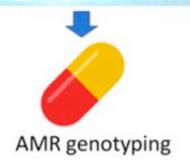
Curated database Web based GUI

Standardized toolkit Analysis workspaces API



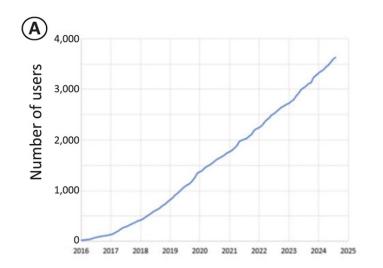


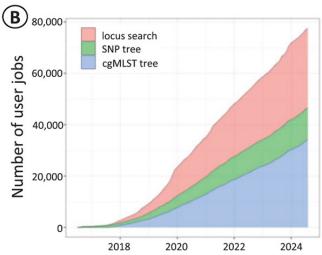




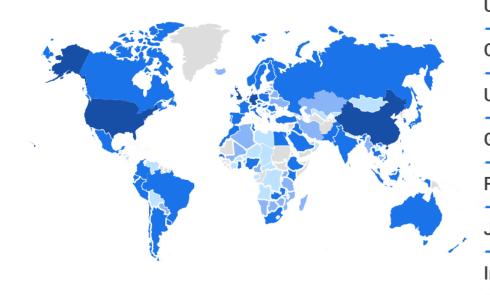
Usage statistics

- (A) Registered EnteroBase users, 2016-2024
- **(B)** Cumulative number of user jobs
- **(C)** User access statistics collected by Google Sep 24 to Sep 25. Table shows top seven countries









Jnited Kingdom	6.3K
China	3.4K
Jnited States	3.1K
Germany	2.2K
rance	1K
Japan	847
ndia	654

~4000 users worldwide

Including Government Agencies:

- US Department of Agriculture
- Public Health England
- Institut Pasteur (France)
- Animal Health and Veterinary Laboratories Agency (UK)
- Antibiotic Resistance Monitoring and Reference Laboratory (UK)
- National Salmonella Reference Library (Ireland)
- Bavarian Health and Food Safety Authority
- Laboratoire National de Sante (Luxembourg)
- National Reference Laboratories (Israel)
- Public Health Agency of Sweden
- Norwegian National Advisory Unit on Detection of Antimicrobial Resistance
- National Institute of Health (Portugal)
- Centre for Zoonoses and Environmental Microbiology (Netherlands)
- Australian National Salmonella Reference Laboratory
- Centre for Infectious Diseases & Microbiology (Sydney, Australia)
- CDC (Taiwan), etc etc etc





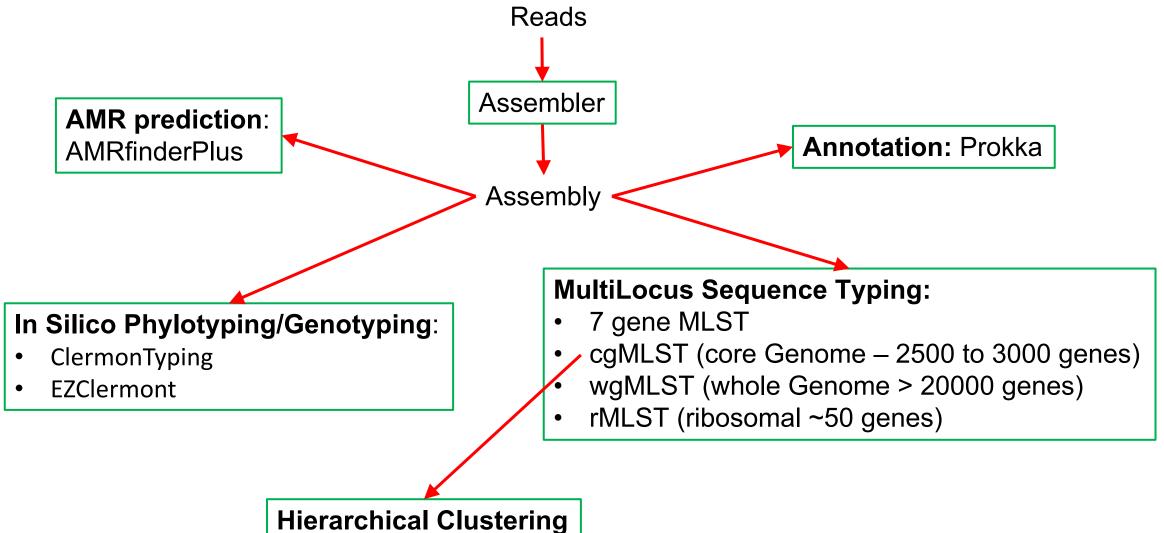






Enterobase data flow and analysis





EToKi Assembler (EnteroBase Tool Kit)



- https://github.com/zheminzhou/EToKi
- SPAdes-based assembly pipeline, with improvements:
- Pre-processing "prepare" step:
 - Trims sequences based on base-qualities (bbduk)
 - Removes potential adapters and barcodes (bbduk)
 - Limits total amount of reads to be used to 120X
- "Assemble" step:
 - Assembles short reads (SPAdes)
 - Maps reads back to assembled genome (Minimap2)
 - Polishes consensus using hapog
 - Removes low level contaminations (contigs with <30% of average read coverage)
 - Estimates the base quality of the consensus
 - Predicts taxonomy using Kraken (MiniKraken)

Assembly tool evaluation exercise for UK government project PATHSAFE



(Pathogen Surveillance in Agriculture, Food and the Environment)

EToKi outperformed:

- SPAdes (default settings)
- MEGAHIT
- IDBA
- Unicycler
- Velvet
- SKESA
- Masurca
- PathogenWatch

Multiple assessment criteria including comparison of post-assembly in silico MLST genotyping

Comparison of hierCC and cgMLST LIN codes



- Strains allocated to clusters using identical algorithm
- Same algorithm for adding strains to existing clusters
- Same algorithm for creating new clusters
- Different ways of identifying the same clusters
- No cluster fusion
- Cluster definitions are stable

Comparison of hierCC and cgMLST LIN codes

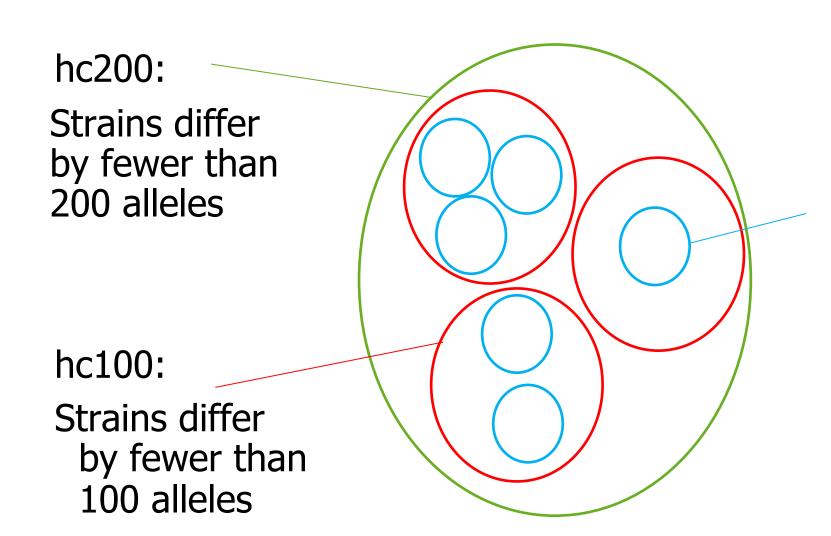


• 1:1 mapping between cluster identifiers.

LIN codes	hierCC
3-0-1-2-3, 3-0-1-2-4	34,56
Explicit cluster hierarchy	Cluster hierarchy not obvious in cluster identifiers
Cluster at level N identified by unique sequence of N numbers	Cluster at any level identified by a unique number. Numbers reused between levels.

Comparison of hierCC and LIN codes



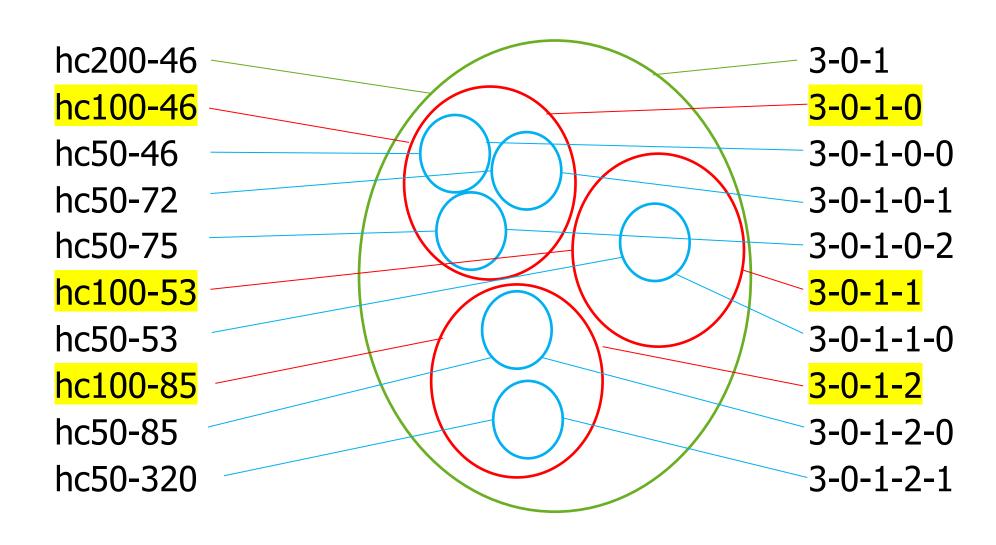


hc50:

Strains differ by fewer than 50 alleles

Comparison of hierCC and LIN codes

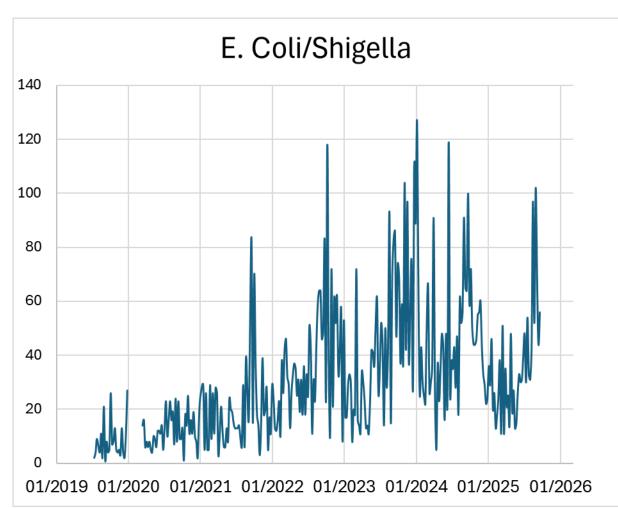


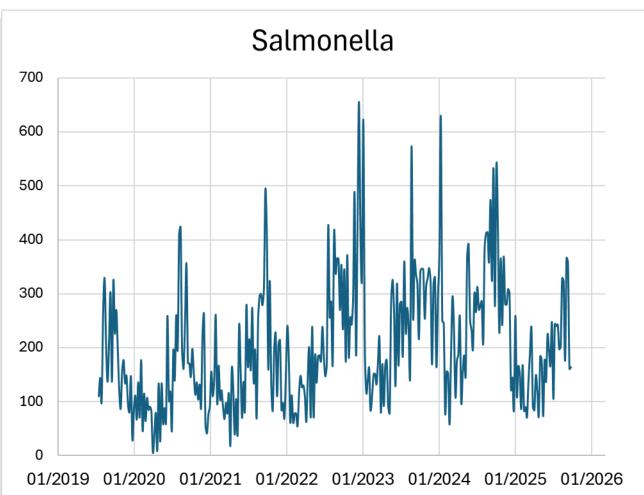


Outbreak surveillance: Institut Pasteur



Weekly strain submission







> 10 000 clinical isolates per year (hospitals and private labs)

Isolation in **XLD** (Xylose Lysine Deoxycholate Agar)

- **Hajna & Chromagar** (in case of contaminations and/or atypical colonies)
- **MALDI-TOF MS** (genus confirmation in case of atypical colonies)
- **Seroagglutination** (in case of several isolates for the same patient)

Whole-genome sequencing (Illumina NextSeq500) & molecular typing:

- Species identification
- MLST (7 locus-scheme, Achtman et al. *PLoS Pathog* 2012)
- resistome/virolome (ResFinder, PlasmidFinder...)
- in silico agglutination patterns (Salmotyper: Fabre, in prep.)
- cgMLST (3,002 genes) -

(Alikhan et al. PLoS Genet 2018)

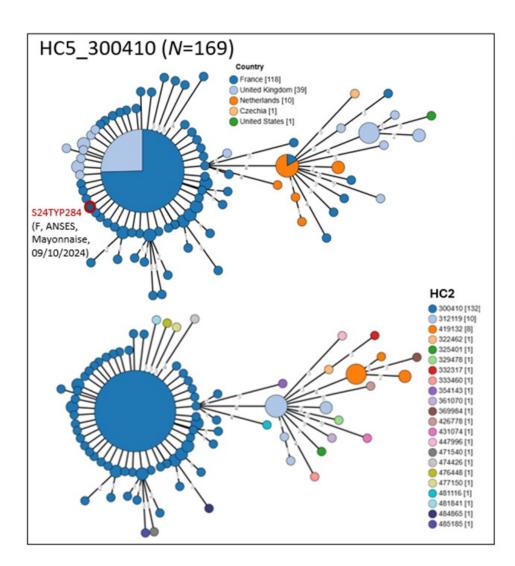
Identification of clusters cgMLST (HC2/HC5 with >8 cas in a 8 weeks timeframe)

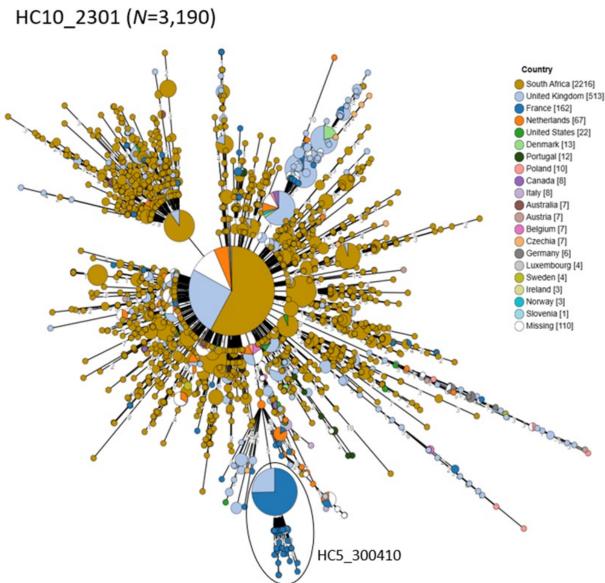
Communication of clusters to **Santé publique France** for investigation in collaboration with food authorities (DGAL, MUS, etc.).

- wgSNP analysis (Enterobase; snippy) confirmation of cgMLST clusters at whole genome level (publications, court investigations, etc.)
- Long-reads sequencing (ONT Nanopore) reconstruction chromosomes and plasmides

Oct-Nov 2024 Salmonella outbreak



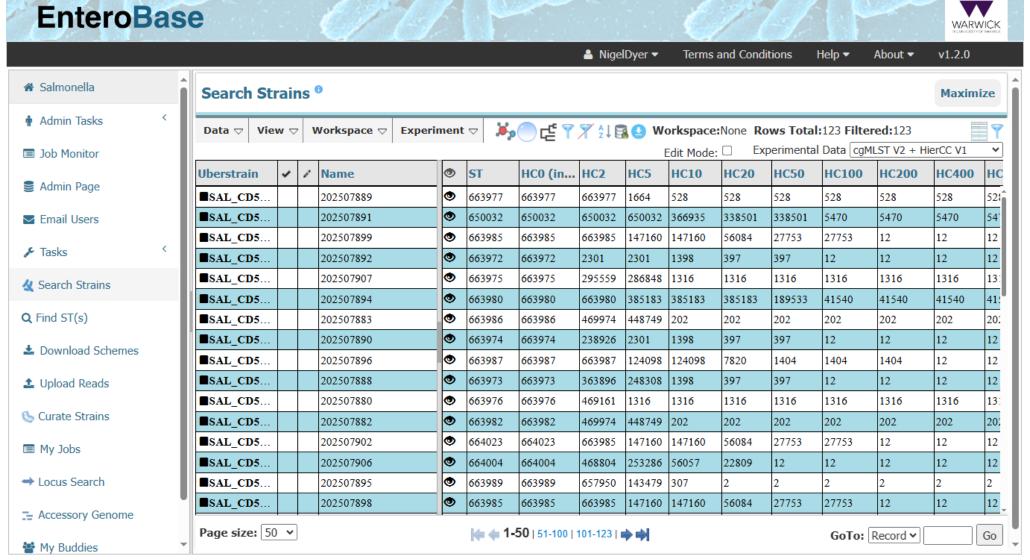




One day's data from France: hierCC identifiers

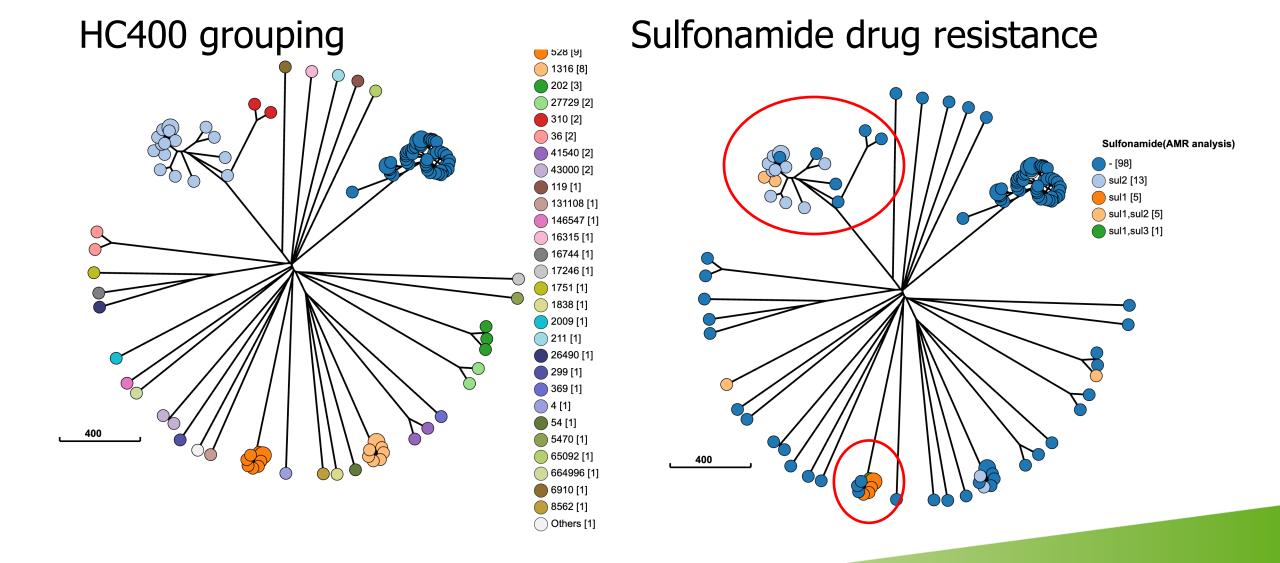






Analysis of one day's data from France

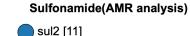


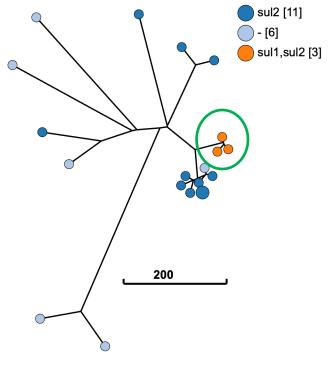


Zoom in on samples of interest



Appears to contain three samples with multi-drug resistance





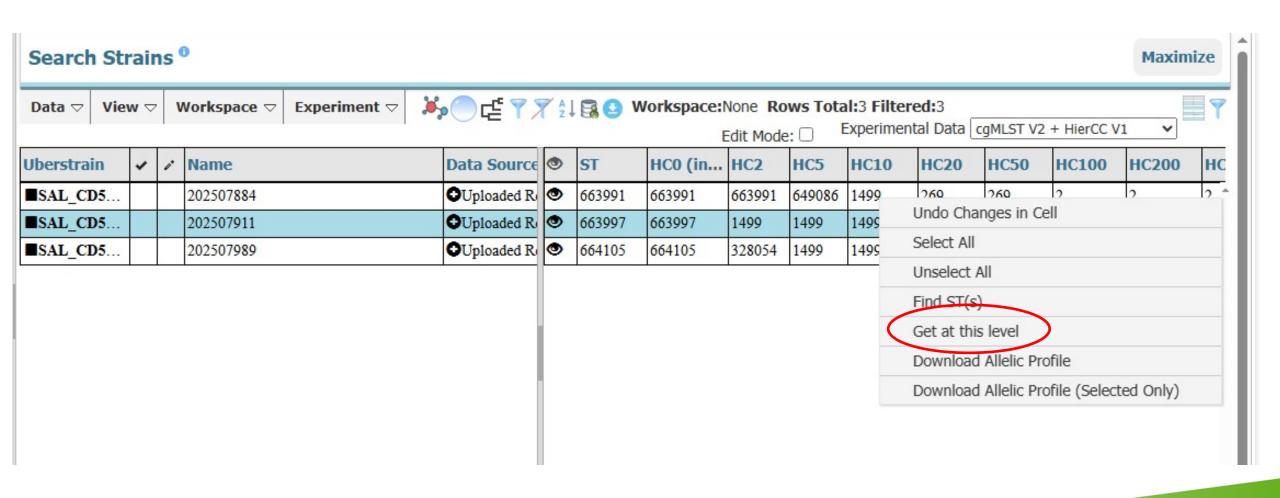
					Edit Mode: Experimental Data						
Uberstrain	~	1	Name		Penicllin	Phenicol	Quinolone	Sulfonamid	Tetracycline	Trimethopr	Other Classes
■SAL_CD5			202507895		blaTEM-1	-	-	sul2	tet(B)	-	-
■SAL_CD5			202507879		blaTEM-1	floR	qnrB19	sul2	tet(A),tet(B)	dfrA12	Lincosamide:lnu(F
SAL_CD5			202507948		blaTEM-1	-	-	sul2	tet(B)	-	-
SAL_CD5			202507951		blaTEM-1	-	-	sul2	tet(B)	-	-
SAL_CD5			202507994		-	-	-	sul2	tet(A)	-	-
SAL_CD5			202507970		blaTEM-1	-	-	sul2	-	-	-
SAL_CD5			202507965		blaTEM-1	-	-	sul2	tet(B)	-	-
SAL_CD5			202507982		-	-	-	sul2	tet(A)	-	-
SAL_CD5			202507973		blaTEM-1	-	-	sul2	-	-	-
■SAL_CD5			202508027		blaTEM-1	floR	-	sul2	tet(B)	-	-
SAL_CD5			202508023	+	blaTEM-1	_	_	su12.	tet(B)	_	_
SAL_CD5			202507884		-	floR	qnrB2	sul1,sul2	tet(A)	dfrA1	-
SAL_CD5			202507911		-	floR	qnrB19	sul1,sul2	tet(A)	dfrA1	-
SAL_CD5			202507989		-	floR	qnrB2	sul1,sul2	tet(A)	dfrA1	-
SAL_CD5			202507908	4	_	-	_		_	_	
SAL_CD5			202507917		-	-	-	-	-	-	-
SAL_CD5			202507995		-	-	-	-	-	-	-
SAL_CD5			202507966		-	-	-	-	-	-	-
SAL_CD5			202507971		-	-	-	-	-	-	-
SAL CD5			202508002								

Click through to associated NCBI data

GenPept → Send to: ▼ Change This record is a non-redundant protein sequence. Please read more here. Custor **MULTISPECIES:** sulfonamide-resistant dihydropteroate **▲** Download Datasets synthase Sul2 [Pseudomonadota] Analyze NCBI Reference Sequence: WP 001043265.1 Run BLA Identical Proteins FASTA Graphics Identify (Go to: (V) Highlight LOCUS 271 aa linear WP 001043265 BCT 03-MAR-2025 Find in the DEFINITION MULTISPECIES: sulfonamide-resistant dihydropteroate synthase Sul2 [Pseudomonadota]. ACCESSION WP 001043265 Protein VERSION WP 001043265.1 KEYWORDS Sulfonar SOURCE Pseudomonadota (proteobacteria) Sul1 Pseudomonadota ORGANISM Total pro Bacteria; Pseudomonadati. Total ger REFSEQ: This record represents a single, non-redundant, protein COMMENT Conserv sequence which may be annotated on many different RefSeq genomes from the same, or different, species. Related ##Evidence-For-Name-Assignment-START## BioProje Evidence Category :: Antimicrobial Resistance Protein Evidence Accession :: NG 048118.1 Nucleotic Evidence Source :: Bacterial Antimicrobial Resistance Reference Gene Database Taxonon ##Evidence-For-Name-Assignment-END## CDD Se COMPLETENESS: full length. FEATURES Location/Qualifiers Conserv source 1..271 /organism="Pseudomonadota" Conserv /db_xref="taxon: 1224" Domain 1..271 gene /gene="sul2" Genome Protein 1...271 /product="sulfonamide-resistant dihydropteroate synthase Genomic Sul2" /calculated_mol_wt=28397 Protein (

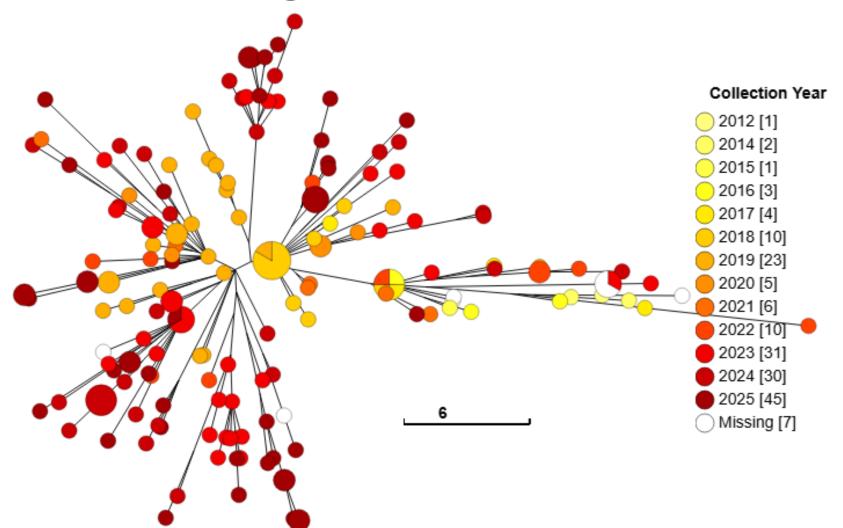
Identify hc10 cluster and find other entries





hc10_499 goes back to 2012

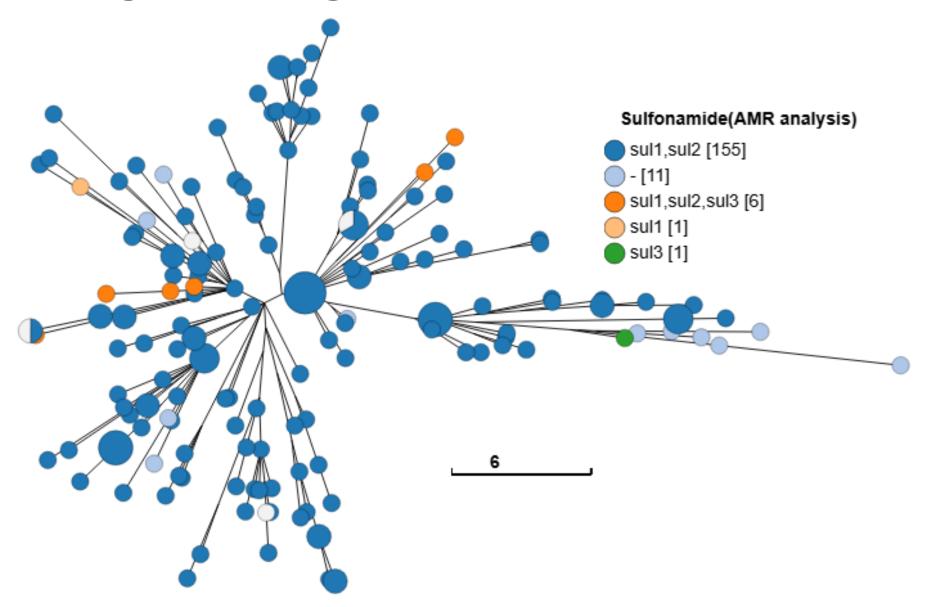




Enterobase holds 178 samples in this cluster, going back to 2012

Long-standing sulfonamide resistance





References



- Zhou Z, Alikhan NF, Mohamed K, the Agama Study Group, Achtman M (2020). The EnteroBase user's guide, with case studies on Salmonella transmissions, Yersinia pestis phylogeny, and Escherichia core genomic diversity, Genome Research, 30:138-152
- Zhou Z, Charlesworth J, Achtman M. (2021). HierCC: A multi-level clustering scheme for population assignments based on core genome MLST, <u>Bioinformatics</u>, <u>37(20):3645–3646</u>
- Nigel P Dyer, Birgitta Päuker, Laura Baxter, Anshul Gupta, Boyke Bunk, Jörg Overmann, Margo Diricks, Viola Dreyer, Stefan Niemann, Kathryn E Holt, Mohammed Rahman, Paul E Brown, Richard Stark, Zhemin Zhou, Sascha Ott, Ulrich Nübel, EnteroBase in 2025: exploring the genomic epidemiology of bacterial pathogens, *Nucleic Acids Research*, 2024, https://doi.org/10.1093/nar/gkae902
- https://enterobase.readthedocs.io/en/latest/

Acknowledgements



Original Developers:

- Mark Achtman
- Zhemin Zhou
- Nabil-Fareed Alikhan
- Martin Sergeant



Acknowledgements



Current Team Members



Laura Baxter

Nigel Dyer

Richard Stark

Paul Brown

Sascha Ott



Birgitta Päuker, Anshul Gupta, Boyke Bunk, Jörg Overmann, Ulrich Nübel

References



- Prokka Annotation: https://github.com/tseemann/prokka
- AMRfinderPlus: https://www.ncbi.nlm.nih.gov/pathogens/antimicrobial-resistance/AMRFinder/
- ClermonTyping: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6113867/
- EZClermont https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7656184/



Acknowledgements

The creation of this training material was commissioned by ECDC to the University of Warwick with the direct involvement of Dr Nigel Dyer